

# 機械学習を用いたゲノムワイド遺伝子多型情報に基づくうつ状態のリスク予測

著者	高橋 雄太
号	89
学位授与機関	Tohoku University
学位授与番号	医博第4035号
URL	<a href="http://hdl.handle.net/10097/00129472">http://hdl.handle.net/10097/00129472</a>

氏名	たかはし ゆうた 高橋 雄太
学位の種類	博士(医学)
学位授与年月日	2020年3月25日
学位授与の条件	学位規則第4条第1項
研究科専攻	東北大学大学院医学系研究科(博士課程) 医科学専攻
学位論文題目	機械学習を用いたゲノムワイド遺伝子多型情報に基づくうつ状態のリスク予測
論文審査委員	主査 教授 富田 博秋 教授 栗山 進一 教授 菅原 準一

## 論文内容要旨

【背景・目的】ゲノムワイド遺伝子多型データを用いてうつ状態のリスクを予測したこれまでの研究では、あまり高い予測精度は示されていない。その原因の一つとして、表現型に実際には効果を持たないバリエーションが多く予測モデルに含まれることにより、予測モデルの学習の段階では見かけ上、高い予測精度が得られるのに、独立したテストデータで検証すると低い予測精度しかえられない、過剰適合と呼ばれる現象が起きてしまうことがある。STMGP (Smooth-Threshold Multivariate Genetic Prediction) 法は過剰適合を軽減させることで予測精度を向上させることを目的に開発された、機械学習を用いた予測モデルである。今回の研究では、ゲノムワイド遺伝子多型データからのうつ状態の予測に STMGP 法を用いて、予測精度を従来法と比較した。

【方法】東北メディカル・メガバンクプロジェクトでリクルートされた、宮城県在住の 3,685 人の遺伝子多型情報を用いて予測モデルを学習させ、岩手県在住の 3,048 人の遺伝子多型情報を用いて、予測モデルの精度を評価した。HumanOmniExpressExome BeadChip Array を用いてゲノタイピングした。うつ症状は Center for Epidemiologic Studies-Depression Scale で評価した。STMGP 法による予測精度と過剰適合の程度は、遺伝子スコア法、GBLUP (genomic best linear unbiased prediction) 法、SBLUP (summary-data-based best linear unbiased prediction) 法、BayesR 法、Ridge 回帰法と比較した。

【結果】STMGP 法による予測精度 (predictive correlation coefficients  $\pm$  標準誤差) は  $0.0769 \pm 0.0173$  であり、遺伝子スコア法 ( $0.0332 \pm 0.0178$ )、GBLUP 法 ( $0.0309 \pm 0.0178$ )、SBLUP 法 ( $0.0164 \pm 0.0178$ )、BayesR 法 ( $0.0100 \pm 0.0185$ )、Ridge 回帰法 ( $0.0260 \pm 0.0178$ ) よりも高かった。また STMGP 法ではトレーニングデータでの見かけ上の予測精度は  $0.3232 \pm 0.0153$  であり、遺伝子スコア法 ( $0.9027 \pm 0.0076$ )、GBLUP 法 ( $0.9627 \pm 0.0017$ )、SBLUP 法 ( $0.9554 \pm 0.0019$ )、BayesR 法 ( $0.9633 \pm 0.0015$ )、Ridge 回帰法 ( $0.9998 \pm 0.0000$ ) よりも低く、過剰適合を軽減させることで予測精度が向上していることが示唆された。

【結論】STMGP 法は過剰適合を軽減することで、従来法よりもゲノムワイド遺伝子多型データからうつ症状を予測する際に、高い精度を示した。微小な効果をもつ遺伝子多型の集合で説明されるような複雑な遺伝疾患の脆弱性の予測に STMGP 法が有効であることが示唆された。

## 審査結果の要旨

博士論文題目 .....機械学習を用いたゲノムワイド遺伝子多型情報に基づくうつ状態のリスク予測.....

所属専攻・分野名 .....医科学 専攻 .....精神神経学 分野.....

学籍番号 B6 MD 5081 氏名 .....高橋 雄太.....

本博士論文の背景として、うつ状態と相関するゲノムワイド遺伝子多型解析は多くなされてきているが、そのデータを用いてうつ状態のリスクを予測する研究では、これまでのところあまり高い予測精度は示されていないことがある。その原因の一つとして、うつ状態という表現型に実際には効果を持たない遺伝子多型を多く予測モデルに含むことにより、予測モデルを形成する学習の段階では見かけ上、高い予測精度が得られるものの、独立したテストデータで検証すると低い予測精度しかえられない、過剰適合と呼ばれる現象が起きてしまうことがあげられる。本論文は、この状況をうけ、過剰適合を軽減させることで予測精度を向上させることを目的に開発された STMGP (Smooth-Threshold Multivariate Genetic Prediction) 法をうつ状態予測の機械学習に導入し、ゲノムワイド遺伝子多型データからのうつ状態の予測精度を従来法と比較したものである。東北メディカル・メガバンクプロジェクトでリクルートされた、宮城県在住の 3,685 人の Human OmniExpress Exome BeadChip Array による遺伝子多型情報を用いて予測モデルを学習させ、岩手県在住の 3,048 人の遺伝子多型情報を用いて、Center for Epidemiologic Studies-Depression Scale で評価したうつ状態の予測モデルの精度を評価した。この STMGP 法による予測精度と過剰適合の程度を、遺伝子スコア法、GBLUP (genomic best linear unbiased prediction) 法、SBLUP (summary-data-based best linear unbiased prediction) 法、BayesR 法、Ridge 回帰法と比較した。STMGP 法による予測精度 (predictive correlation coefficients  $\pm$  標準誤差) は  $0.0769 \pm 0.0173$  であり、遺伝子スコア法 ( $0.0332 \pm 0.0178$ )、GBLUP 法 ( $0.0309 \pm 0.0178$ )、SBLUP 法 ( $0.0164 \pm 0.0178$ )、BayesR 法 ( $0.0100 \pm 0.0185$ )、Ridge 回帰法 ( $0.0260 \pm 0.0178$ ) よりも高かった。また STMGP 法ではトレーニングデータでの見かけ上の予測精度は  $0.3232 \pm 0.0153$  であり、遺伝子スコア法 ( $0.9027 \pm 0.0076$ )、GBLUP 法 ( $0.9627 \pm 0.0017$ )、SBLUP 法 ( $0.9554 \pm 0.0019$ )、BayesR 法 ( $0.9633 \pm 0.0015$ )、Ridge 回帰法 ( $0.9998 \pm 0.0000$ ) よりも低く、過剰適合を軽減させることで予測精度が向上していることが示唆された。STMGP 法は過剰適合を軽減することで、従来法よりもゲノムワイド遺伝子多型データからうつ症状を予測する際に、高い精度を示すことが証明された。本研究はゲノム情報に基づくうつ状態の予測精度の向上を行なったのみならず、微小な効果をもつ遺伝子多型の集合で説明されるような複雑な遺伝疾患の脆弱性の予測に STMGP 法が有効であることを示唆したことでも意義が認められる。よって、本論文は博士（医学）の学位論文として合格と認める。